

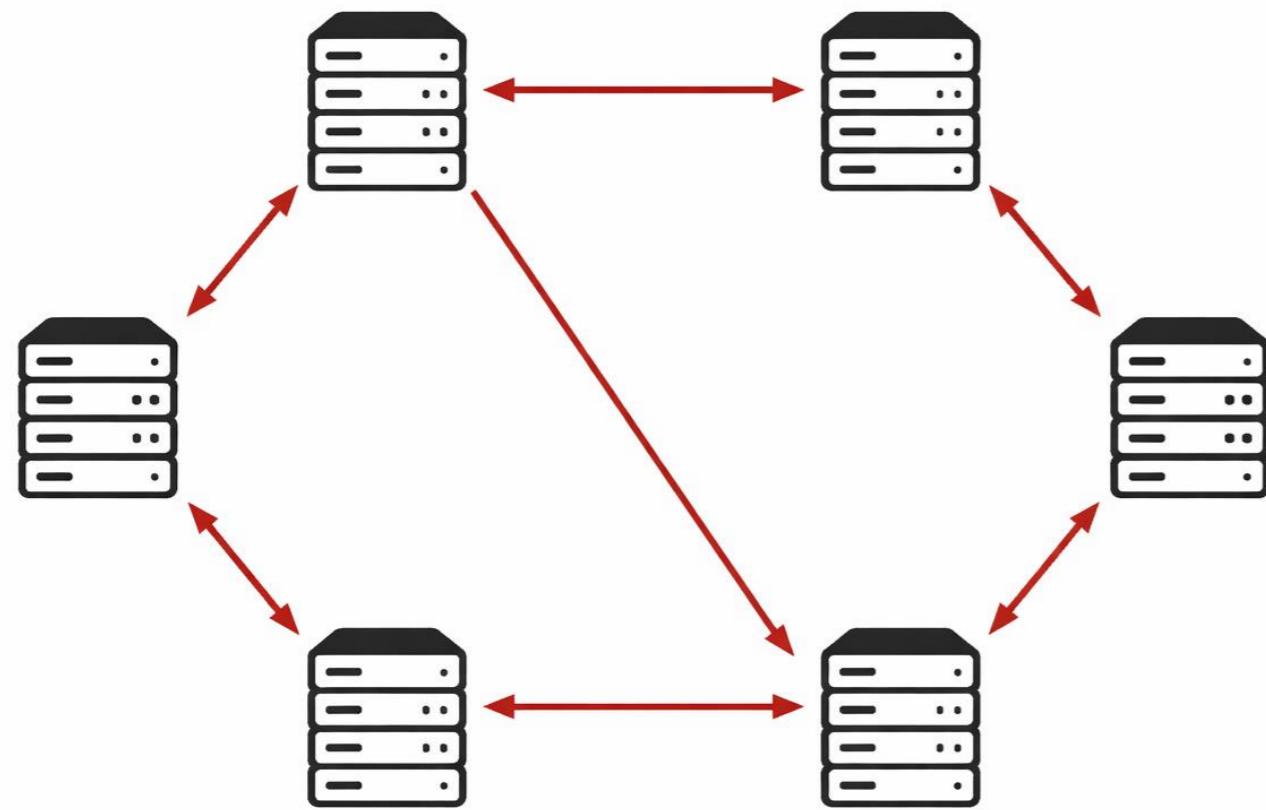
# **Accelerated Dual Method for Distributed Optimization: An Inexact-Gradient View of Local Updates**

Junchi Yang, Ziyang Zeng, Linxuan Pan, Murat Yildirim, Feng Qiu

Presenter: Junchi Yang (Chinese University of Hong Kong, Shenzhen)

**ICML 2026**

# Distributed Optimization



$$\min_{x \in \mathbb{R}^n} F(x) := \frac{1}{M} \sum_{i=1}^M f_i(x).$$

# Local Updates for Distributed Optimization

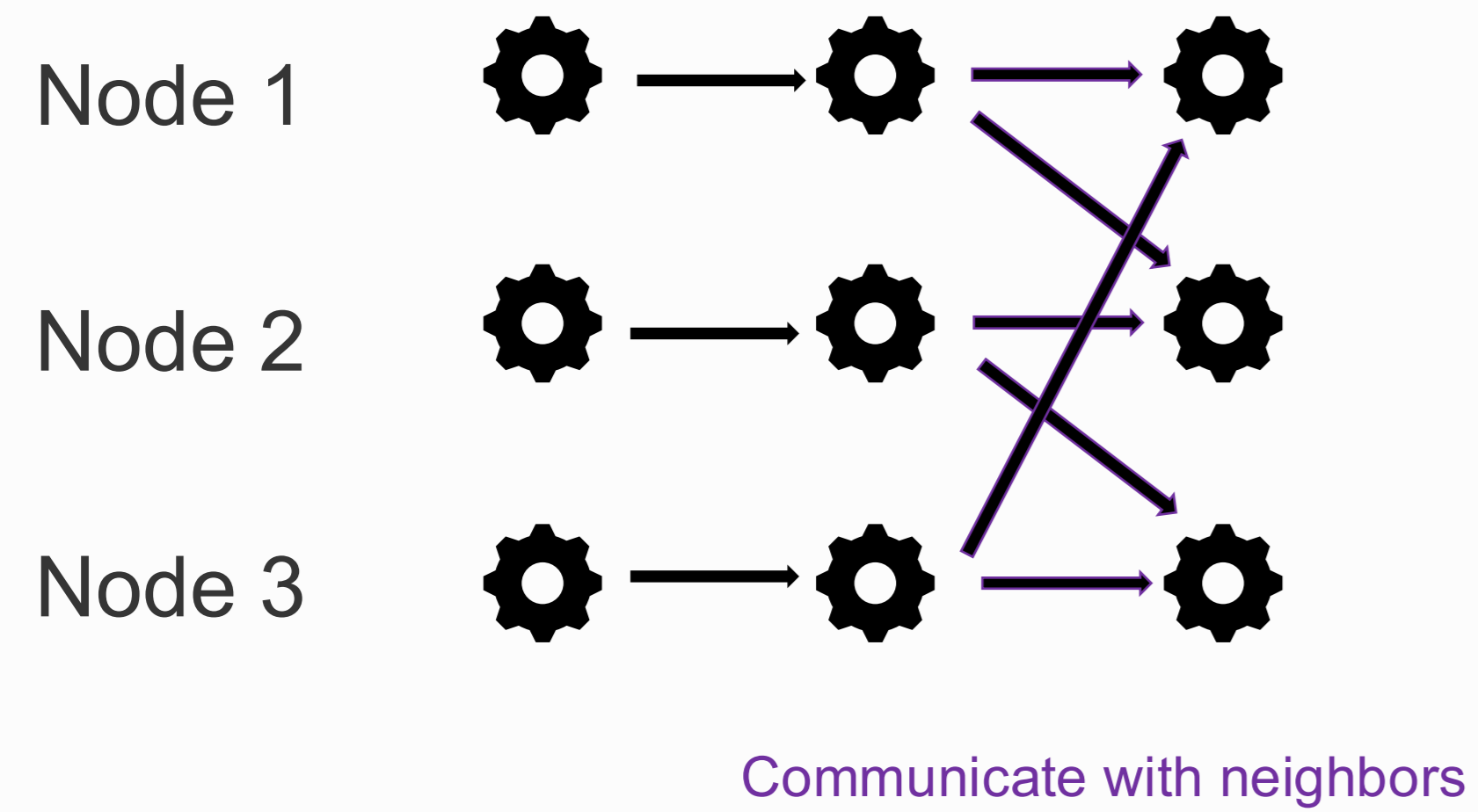
## Without Local Updates



Compute one local gradient/  
One local update

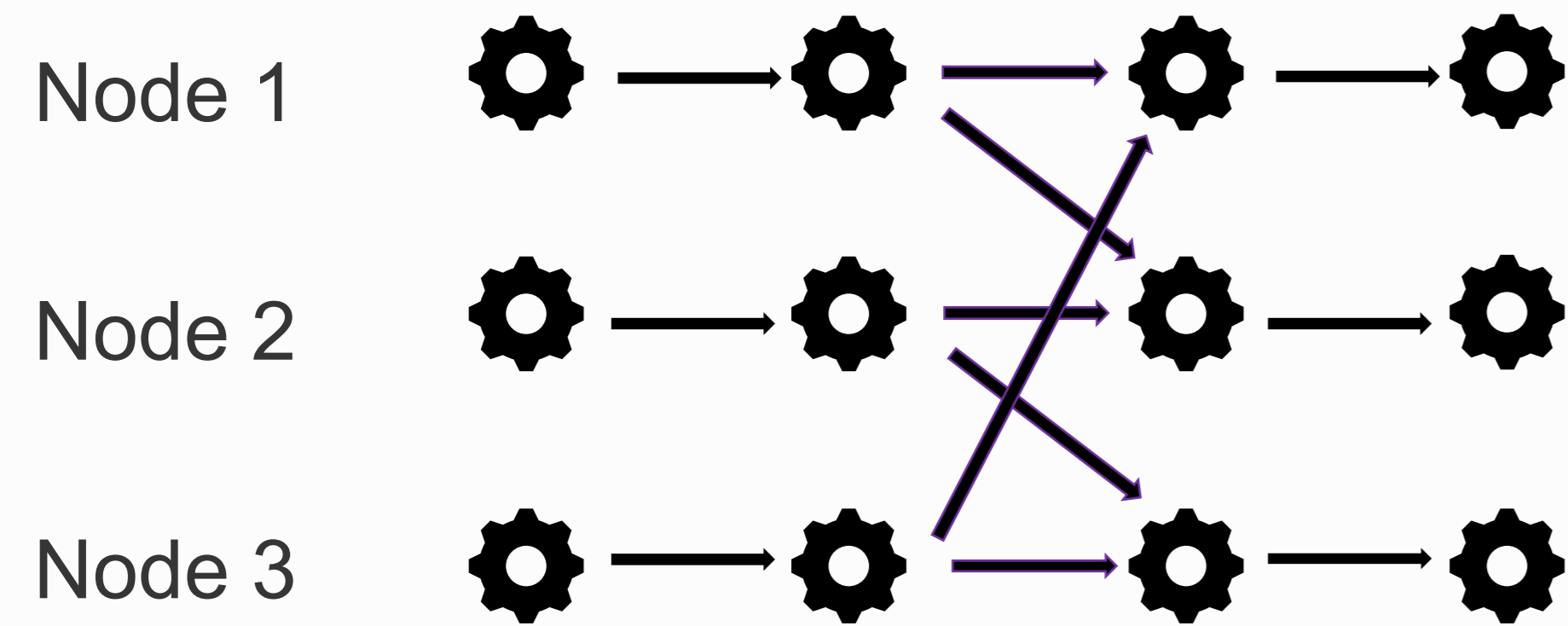
# Local Updates for Distributed Optimization

Without Local Updates



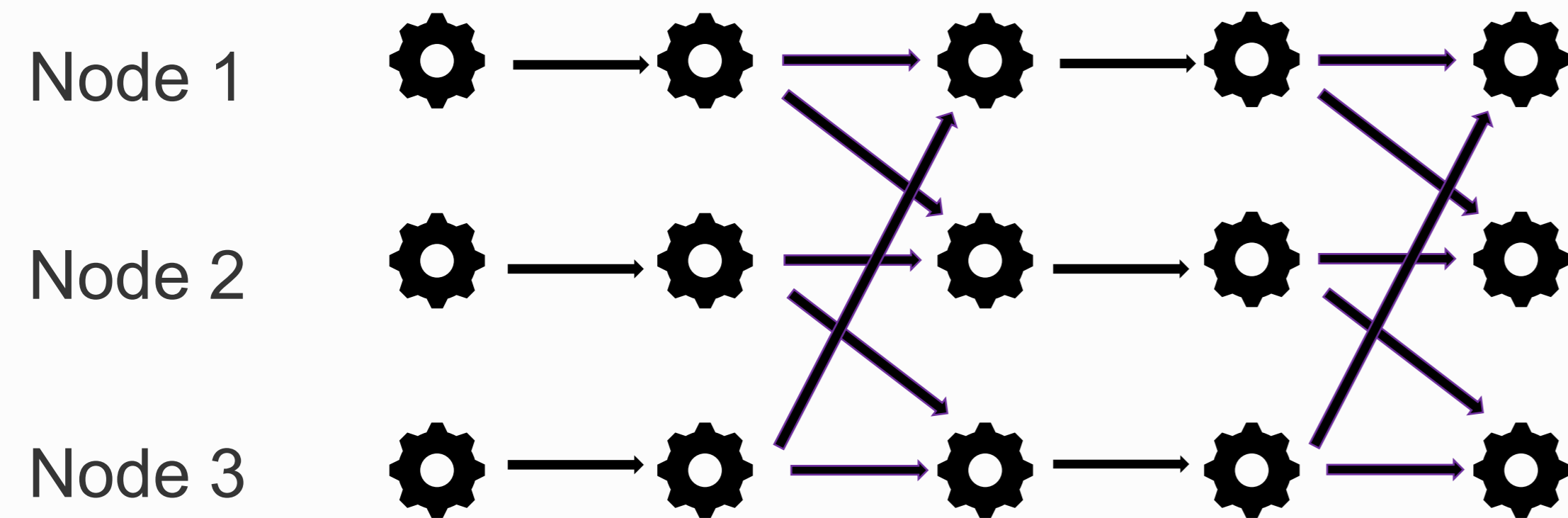
# Local Updates for Distributed Optimization

Without Local Updates



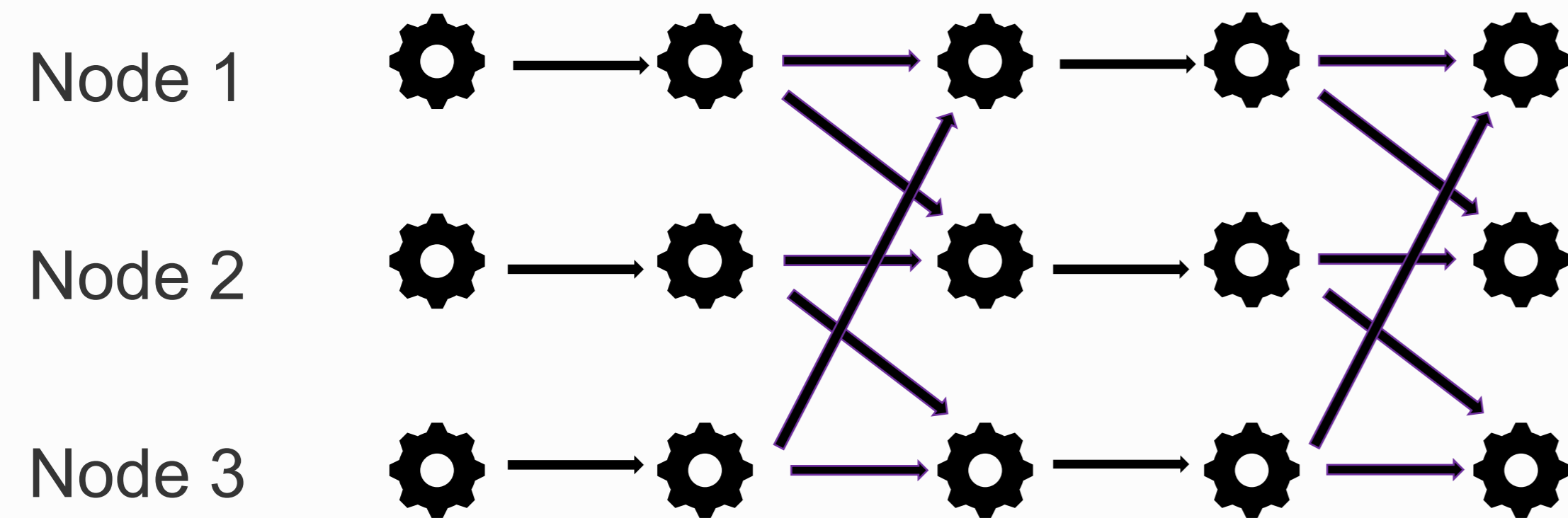
# Local Updates for Distributed Optimization

Without Local Updates

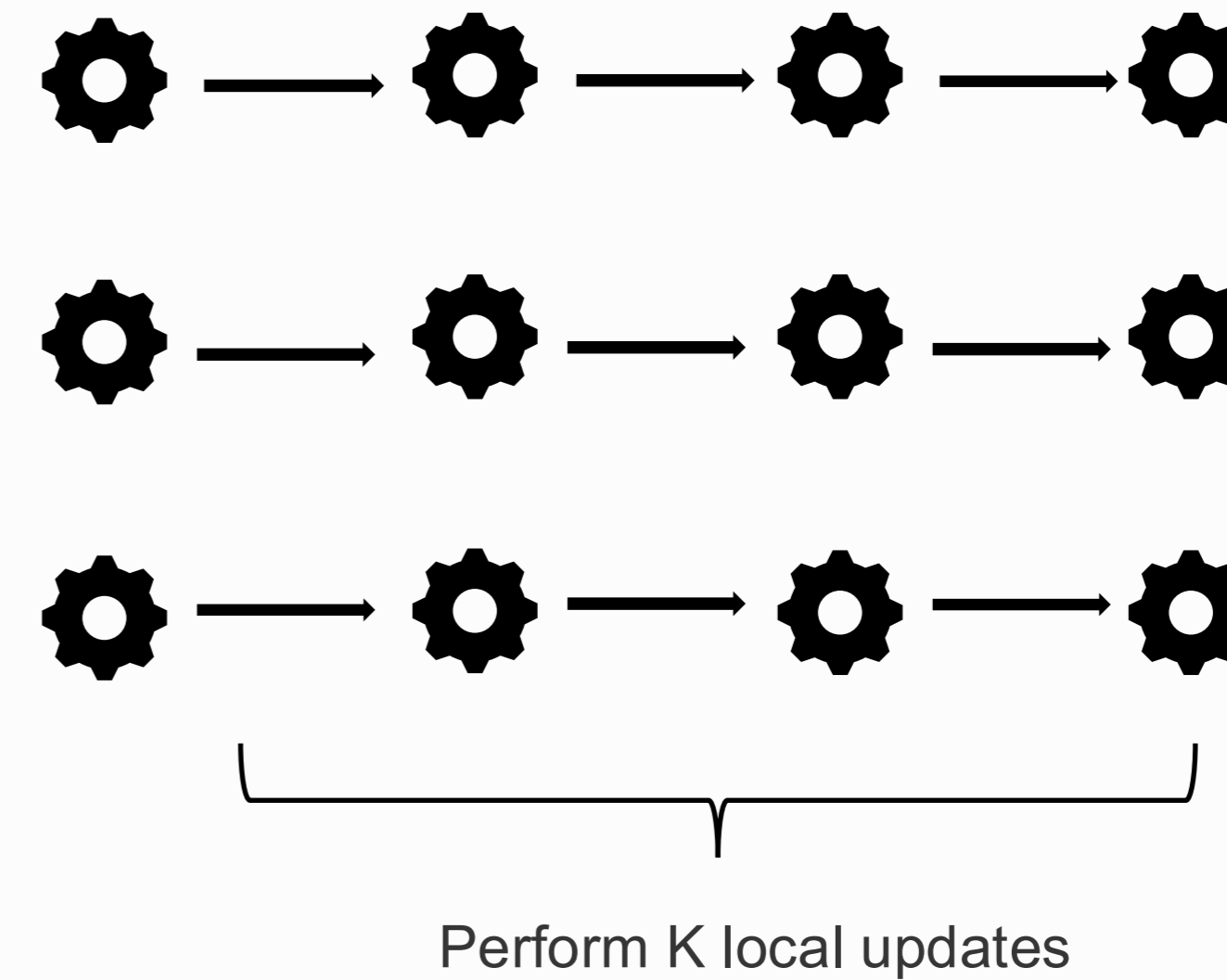


# Local Updates for Distributed Optimization

## Without Local Updates

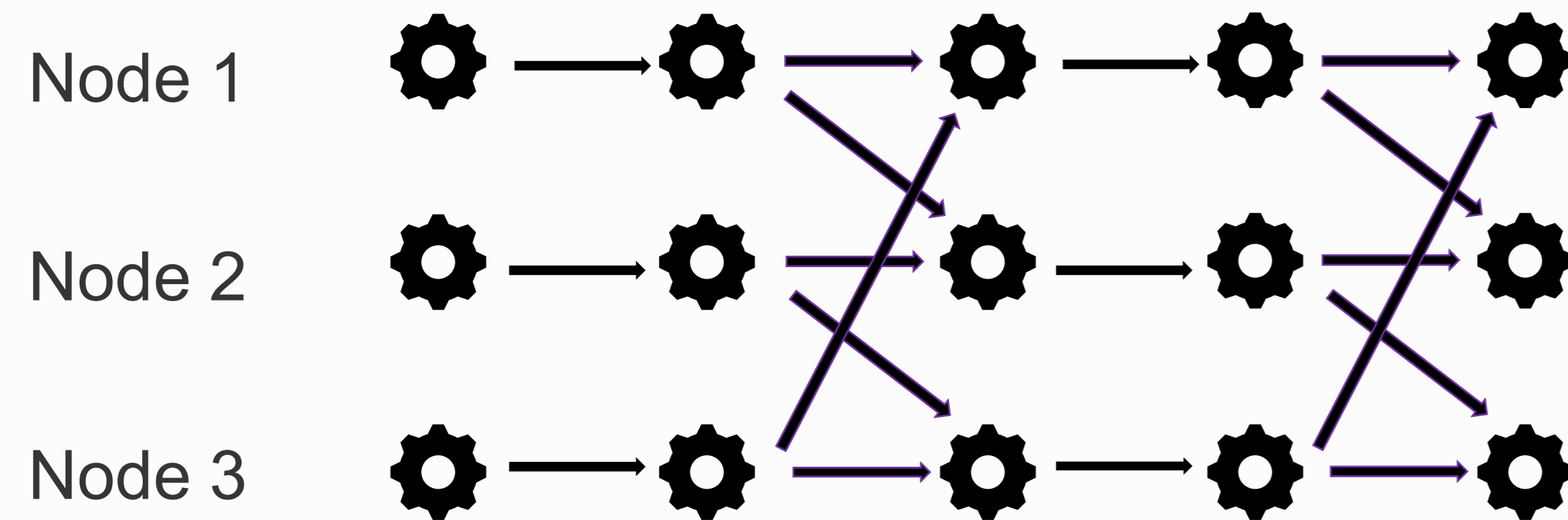


## With Local Updates

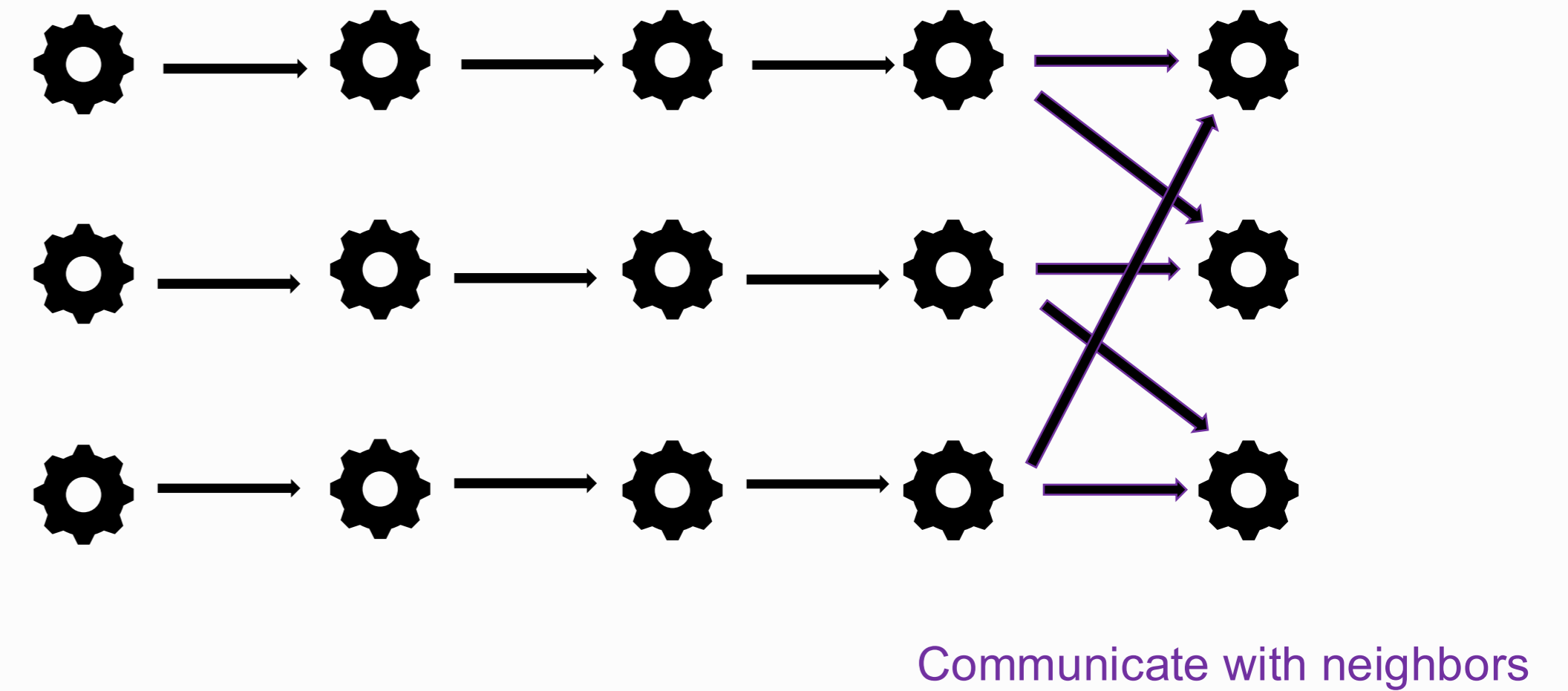


# Local Updates for Distributed Optimization

## Without Local Updates



## With Local Updates



# Strongly Convex Stochastic Distributed Optimization

- ▶ **Setting:**  $\mu$ -strong convexity +  $\ell$ -Lipschitz smoothness

Condition number  $\kappa = \ell/\mu$ .

- ▶ **The mixing matrix  $W$ :** doubly stochastic, symmetric, and graph connected.

Denote  $p = 1 - \sigma_2$ , where  $\sigma_2$  is the second largest singular value of  $W$ .

# Accelerated Dual Methods

- ▶ (Accelerated) dual methods have been widely studied [Terelius et al., 2011; Ghadimi et al., 2011; Scaman et al., 2017; Uribe et al., 2020].
- ▶ Define  $U = (I - W)^{\frac{1}{2}}$ . The original problem can be reformulated as

$$\begin{aligned} \min_{X \in \mathbb{R}^{M \times n}} \quad & H(X) && \text{where } X = (x^1, \dots, x^M)^\top \in \mathbb{R}^{M \times n} \\ \text{s.t.} \quad & UX = 0 && H(X) = \frac{1}{M} \sum_{i=1}^M f_i(x^i) \end{aligned}$$

- ▶ Lagrangian for this constrained problem:

$$\min_{X \in \mathbb{R}^{M \times n}} \max_{\lambda \in \mathbb{R}^{M \times n}} \mathcal{L}(X, \lambda) = H(X) + \langle \lambda, UX \rangle.$$

# Accelerated Dual Methods

- ▶ Apply  $K$  steps of stochastic gradient descent on  $X$  and accelerated gradient ascent in  $\lambda$

**Primal:** for  $k = 0, \dots, K - 1$

$$X_{t,k+1} = X_{t,k} - \tau_2 (\tilde{\nabla} H(X_{t,k}) + U \tilde{\lambda}_t)$$

**Dual:**  $\lambda_{t+1} = \tilde{\lambda}_t + \tau_1 U X_{t,K}$

$$\tilde{\lambda}_{t+1} = \lambda_{t+1} + \beta (\lambda_{t+1} - \lambda_t)$$

No communication  
 $K$  Local updates

- ▶ Change of variable  $\zeta = U\lambda$ ,

**Primal:** for  $k = 0, \dots, K - 1$

$$X_{t,k+1} = X_{t,k} - \tau_2 (\tilde{\nabla} H(X_{t,k}) + \tilde{\zeta}_t)$$

**Dual:**  $\zeta_{t+1} = \tilde{\zeta}_t + \tau_1 W X_{t,K}$

$$\tilde{\zeta}_{t+1} = \zeta_{t+1} + \beta (\zeta_{t+1} - \zeta_t)$$

**Primal (agent  $i$ ):** for  $k = 0, \dots, K - 1$

$$x_{i,t,k+1} = x_{i,t,k} - \tau_2 (\tilde{\nabla} f_i(x_{i,t,k}) + \tilde{\zeta}_{i,t})$$

**Dual (agent  $i$ ):**  $\zeta_{i,t+1} = \tilde{\zeta}_{i,t} + \tau_1 \sum_{j=1}^M w_{ij} x_{j,t,K}$

$$\tilde{\zeta}_{i,t+1} = \zeta_{i,t+1} + \beta (\zeta_{i,t+1} - \zeta_{i,t})$$

# Accelerated Dual Methods

## Proposition

The dual function  $\Psi(\lambda) = \min_X \mathcal{L}(X, \lambda)$  is  $\mu_\Psi$ -strongly concave and  $\ell_\Psi$ -smooth in  $\text{Span}(U)$ , where  $\mu_\Psi = Mp/\ell$  and  $\ell_\Psi = 2M/\mu$ .

► The condition number of  $\Psi$  is  $2\kappa/p$ . Applying Accelerated Gradient Ascent (AGA) to  $\Psi$  has the communication complexity of  $\tilde{O}\left(\sqrt{\frac{\kappa}{p}}\right)$ .

► This optimal cor



Does it converge under arbitrary  $K$ ?

► When  $K$  is large



Can it attain optimal communication complexity with  $K = O(\epsilon^{-1})$ ?

# Convergence of Accelerated Dual Methods

## Theorem (Accelerated Dual Method)

With proper stepsize and momentum parameter, if the local steps

1.  $K \leq \tilde{\Theta} \left( \max \left\{ \frac{c_0 \sigma^2}{\epsilon}, \kappa \right\} \right)$ , then the number of communication rounds to achieve an  $\epsilon$ -accurate solution is  $\tilde{O} \left( \frac{c_1}{\epsilon K} \right)$ .
2.  $K \geq \tilde{\Theta} \left( \max \left\{ \frac{c_0 \sigma^2}{\epsilon}, \kappa \right\} \right)$ , then the number of communication rounds to achieve an  $\epsilon$ -accurate solution is  $\tilde{O} \left( \left( \frac{\kappa}{p} \right)^\alpha \right)$  with  $\alpha \in \left[ \frac{1}{2}, 1 \right]$ , given  $\beta = 1 - O((\mu\tau_1)^{1-\alpha})$ .

► When local updates  $K \leq \tilde{\Theta} \left( \max \left\{ \frac{c_0 \sigma^2}{\epsilon}, \kappa \right\} \right)$ , the total sample complexity is  $\tilde{O} \left( \frac{c_1}{\epsilon} \right)$ .

► When local updates  $K \geq \tilde{\Theta} \left( \max \left\{ \frac{c_0 \sigma^2}{\epsilon}, \kappa \right\} \right)$ , the communication complexity can be  $\tilde{O} \left( \sqrt{\frac{\kappa}{p}} \right)$ .

# Summary for Strongly Convex Setting

Communication Complexities  
when the number local steps/minibatch is large enough

	Communication Complexity	Local Updates	Arbitrary $K$
<b>Local-DSGD</b> (Koloskova et al., 2020)	$\tilde{O}\left(\kappa p^{-1} + \sqrt{\kappa} p^{-1} \epsilon^{-\frac{1}{2}}\right)$	Yes	Yes
<b>LED</b> (Alghunaim, 2024)	$\tilde{O}(\kappa^2 p^{-1})$	Yes	Yes
<b>Stochastic GT</b> (Koloskova et al., 2021)	$\tilde{O}(\kappa p^{-1} c^{-1})$	No	Yes
<b>Distributed FGM</b> (Uribe et al., 2020)	$\tilde{O}\left(\kappa^{\frac{1}{2}} p^{-\frac{1}{2}}\right)$	Yes	No
<b>Local ADA</b> (our work)	$\tilde{O}\left(\kappa^{\frac{1}{2}} p^{-\frac{1}{2}}\right)$	Yes	Yes
<b>Lower Bound</b> (Scaman et al., 2017)	$\tilde{\Omega}\left(\kappa^{\frac{1}{2}} p^{-\frac{1}{2}}\right)$	N/A	N/A